

OPTIMIZING HUMAN-INTERPRETABLE DIALOG MANAGEMENT POLICY USING GENETIC ALGORITHM

Hang Ren, Weiqun Xu and Yonghong Yan

The Key Laboratory of Speech Acoustics and Content Understanding
Institute of Acoustics, Chinese Academy of Sciences
21 North 4th Ring West Road, Beijing, China, 100190

ABSTRACT

Automatic optimization of spoken dialog management policies that are robust to environmental noise has long been the goal for both academia and industry. Approaches based on reinforcement learning have been proved to be effective. However, the numerical representation of dialog policy is human-incomprehensible and difficult for dialog system designers to verify or modify, which limits its practical application. In this paper we propose a novel framework for optimizing dialog policies specified in domain language using genetic algorithm. The human-interpretable representation of policy makes the method suitable for practical employment. We present learning algorithms using user simulation and real human-machine dialogs respectively. Empirical experimental results are given to show the effectiveness of the proposed approach.

Index Terms— dialog management, reinforcement learning, genetic algorithm

1. INTRODUCTION

Dialog manager (DM) plays a central part in spoken dialog system (SDS) and its major functionalities include tracking dialog states and maintaining a dialog policy which decides how the system reacts given certain dialog state. Designing a dialog policy by hand is tedious and erroneous because of the uncertainty of underlying dialog states especially in noisy environment. In recent years various approaches for automatic DM policy optimization have been proposed [1, 2, 3, 4], among which methods based on reinforcement learning (RL) and POMDP model are the most popular [5]. The main objective of RL is to learn an optimum policy conducted by an agent by maximizing its cumulative reward. One of the advantages of RL-based DMs is its robustness to noises from automatic speech recognizer (ASR) and spoken language understanding (SLU) modules. Also, it automates the optimization process by allowing the agent to discover the optimum policy through exploring the underlying state-action space and incrementally improve the controlling policy.

Despite all the advantages, RL-based DMs are not widely deployed for commercial SDSs due to several reasons [6]. Firstly, RL algorithms are mostly data-demanding, which leaves dialog system designers in a dilemma since there is usually few or even no data available at the early stage of system development. Several methods have been proposed to mitigate this problem. A user simulator is often firstly built using wizard-of-oz dialog data, and then the simulator is used to train a RL-based DM. In recent studies it has been shown that by incorporating domain knowledge into the design of kernel functions, the GPSARSA [7, 8] algorithm exhibits much faster learning speed than conventional online RL methods. Secondly, RL algorithms usually use complex numerical models in optimizing the value function, which are usually beyond human comprehension. The learned policy is implicitly represented in the optimized value function (Q-function), which is difficult or even impossible for system designer to verify or modify, keeping back domain experts from setting necessary constraints over the system behavior.

In this paper we propose to use Genetic Algorithm (GA) [9] in optimizing DM policies (GA-DM) which are comprehensible to human designers and easy to verify and modify. The underlying idea is intuitive. We use human-readable domain language to sketch the basic structure of the DM policy, and leave the uncertain parameters for later tuning. According to our experiences in deploying SDSs, it is relatively easy to specify a basic DM policy, when engineering slot-filling or task-driven SDSs of a moderate scale. The most difficult part lies in setting various threshold parameters in dealing with ASR and SLU errors via repeatedly confirming and grounding. These parameters are usually set heuristically or by trial-and-error. Automatic optimization of these parameters will be of great help. We hope to keep the trade-off between purely hand-designed rule-based policies and the ones automatically learned using black-box and data-driven RL methods while keeping the merits from both approaches. Two variants of the approach are proposed and evaluated, an on-line training method through interaction with a simulated user and an off-line and sample-efficient version called on-corpus Q-points regression.

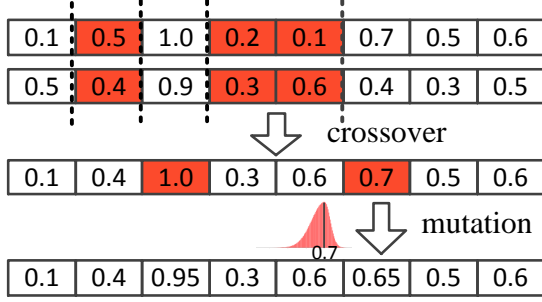


Fig. 1: Crossover and mutation of individuals (chromosomes). Each individual is a real vector with constituent scalars in $[0, 1]$.

Listing 1 BNF grammar of dialog policy template

```

<template> ::= 'if' <cond-exp> 'then' <action> 'else'
              <template>
              | 'if' <cond-exp> 'then' <action> 'else' <action>

<cond-exp> ::= <cond-exp> <logic-op> <cond-exp>
              | <boolean-state-var>
              | <num-state-var> <comparator> <free-param>

<comparator> ::= '<' | '>' | '=='

<logic-op> ::= 'and' | 'or'

```

In the following sections we describe the algorithms and experiments in detail. In section 2.1 we briefly describe Genetic Algorithm and its application in DM policy optimization. We propose two different policy optimization methods based on simulation and dialog corpus in sections 2.2 and 2.3 respectively. In section 3 we give experimental results on simulated user and real human-machine dialog corpus.

2. MODELS AND ALGORITHMS

2.1. Genetic algorithm and dialog policy template

Genetic algorithm is a general optimization framework. It simulates the evolution process of natural selection by keeping a *population* of candidate solutions (individuals) and incrementally improve the quality using various genetic operators. It is a *global* optimization method which can solve both numerical and combinatorial problems. The key constituent of GA is a *fitness* function evaluating the utility of each individual. GA has been proved to be effective in solving various problems, including optimizing controllers in AI games. The pseudocode of optimizing DM using GA is given in Algorithm 1. We refer readers to [9] for a detailed description of GA. The concepts of *genotype* and *phenotype* are not discriminated here. In GA an individual directly carries all the information comprising a solution, which is a fixed-length

Algorithm 1 Genetic algorithm policy optimization

```

1: Input fitness function  $F$ ,  $N_{pop}$ ,  $N_{mut}$ ,  $T_{max}$ ,  $K$ ,  $\sigma$ ,  $\mu_{mut}$ 
2:  $t \leftarrow 0$ ,  $P_0 \leftarrow \emptyset$   $\triangleright$  the initial population
3: for  $i \leftarrow 1, \dots, N_{pop}$  do
4:    $P_0.add(Random.generateIndividual())$   $\triangleright$  random initialization
5:  $P_0.evalFitness()$   $\triangleright$  evaluate fitness of each individual
6: while fitness  $f_t$  not converges and  $t < T_{max}$  do
7:    $t \leftarrow t + 1$ ,  $P_t \leftarrow \emptyset$   $\triangleright$  next generation
8:    $P_t.add(P_{t-1}.getFittest())$   $\triangleright$  elitism
9:   for  $i \leftarrow 1, \dots, N_{mut}$  do
10:     $P_t.add(mutate(P_{t-1}.getFittest(), \sigma, \mu_{mut}))$   $\triangleright$  mutate the fittest
11:   for  $i \leftarrow 1, \dots, N_{pop} - N_{mut} - 1$  do
12:     $I_1, I_2 \leftarrow tournamentSelect(P_{t-1}, K)$ 
13:     $P_t.add(mutate(crossover(I_1, I_2), \sigma, \mu_{mut}))$   $\triangleright$  reproduction
14:    $P_t.evalFitness()$ 
15:    $f_t = P_t.getFittest().getFitness()$ 
16: return  $P_t.getFittest()$ 

17: function MUTATE( $I$ ,  $\sigma$ ,  $\mu_{mut}$ )  $\triangleright$  mutate an individual  $I$ 
18:   for each parameter  $\theta_i$  of  $I$  do
19:     if  $Random.uniform() < \mu_{mut}$  then
20:        $I.\theta_i \leftarrow perturb(I.\theta_i, \sigma)$ 
21:   return  $I$ 

22: function PERTURB( $\theta$ ,  $\sigma$ )  $\triangleright$  add random noise to a single parameter
23:    $g \leftarrow abs(Random.stdGaussian())$ 
24:   if  $Random.uniform() < \theta$  then
25:      $v \leftarrow -\frac{g}{\sigma} * \theta + \theta$ 
26:   else
27:      $v \leftarrow \frac{g}{\sigma} * (1.0 - \theta) + \theta$ 
28:   if  $v < 0.0$  or  $v > 1.0$  then
29:     return  $perturb(\theta, \sigma)$ 
30:   else
31:     return  $v$ 

32: function TOURNAMENTSELECT( $P$ ,  $K$ )  $\triangleright$  tournament selection
33:   choose a random subset  $P_K$  of size  $K$  from  $P$ 
34:   return  $P_K.getFittest()$ 

35: function CROSSOVER( $I_1$ ,  $I_2$ )  $\triangleright$  crossover of two parents
36:    $I' \leftarrow$  exchange random parts of  $I_1$  and  $I_2$ 
37:   return  $I'$ 

```

floating-point array in our experiment and each number is in $[0, 1]$ as a free parameter of the dialog policy template. An individual can instantiate a concrete DM policy, with a defined policy template. The policy template is composed of a set of prioritized condition-action expressions and used to specify the basic structure of a dialog policy. Given certain dialog

state, each condition expression is checked sequentially and the first matched one is selected with the associated action chosen as output. Listing 1 gives the BNF grammar of the proposed templates. The actions of the template is fixed and free parameters can be used to set thresholds for numerical state variables. Apart from the conditional expression, parameters can also be used to induce new state variables, for example a variable representing the number of slots whose top scores are above certain threshold. Although the general system action is fixed in the template, the ‘structure’ of the action (in this slot-filling setting, structure includes sub-dialog-actions and the associated slots and values) can be controlled by parameters. For example, in the action ‘offer’, threshold can be used to filter the hypotheses that are used in searching for the queried information.

Note that the template in Listing 1 has been proposed for its conciseness and simplicity and does not have to be fixed. The design of the dialog template requires knowledge in the dialog domain but does not need a exact model of the environmental noise, thus is very suitable for human experts. This engineering division is intentionally made in our proposed approach.

In GA two kinds of genetic operators are used, i.e. *mutation* and *crossover*, which are shown graphically in Figure 1 and as pseudo-code in Algorithm 1. During crossover, two parents are selected, then random parts of the two parents are exchanged, giving birth to a new child. The mutation operator checks each component of a chromosome sequentially, either leaving it intact or perturbing it randomly. In our implementation the perturbation is realized by sampling from a skewed normal distribution with the peak centered at the perturbed number. If the sampling result lies outside $[0, 1]$, the process is repeated by calling the function *perturb* recursively. This sampling sub-routine is designed for a smooth distribution function. The *mutation* and *crossover* operators represent asexual and bisexual reproductions in GA respectively. Other reproducing strategy can be used as long as it effectively explores the underlying solution space. *Tournament selection* is used to select individuals for reproduction. It is a simple selection method where random K individuals are chosen from the population. We also use the *elitism* technique passing the fittest individual directly to the next generation, ensuring that the fitness of the population will never decrease. The fitness function is the most important part of GA since it guides the algorithm in searching for optimum solution. Two kinds of DM policy fitness evaluation methods are described in the following sections.

2.2. Policy optimization with a user simulator

Since the fitness function should be consistent with the performance of the DM, the most straightforward way is to evaluate it online with users. But interacting with real user is time-consuming and labor-intensive, thus an agenda-based user

Algorithm 2 Episodic fitted Q-iteration

```

1: Input  $\{(s_{i,t}, a_{i,t+1}, s_{i,t+1})\}$  where  $t \leftarrow 1, \dots, T_t - 1$ ,
   and  $i \leftarrow 1, \dots, N$ 
2: initialize Q-function approximator  $\hat{Q}(s, a)$  and array  $Q_{i,t}$ 
   to 0
3: for  $l \leftarrow 1, \dots, L_{max}$  do
4:   for  $i \leftarrow 1, \dots, N$  do ▷ for each dialog
5:     for  $t \leftarrow 1, \dots, T_i - 1$  do ▷ for each turn
6:        $r \leftarrow \text{reward}(s_{i,t}, a_{i,t+1}, s_{i,t+1})$ 
7:       if  $t == T_i - 1$  then
8:          $Q_{i,t} \leftarrow r$  ▷ when the dialog ends
9:       else
10:         $Q_{i,t} \leftarrow r + \gamma \max_a \hat{Q}(s_{i,t+1}, a)$ 
11:      Regress  $\hat{Q}(s, a)$  on  $\{(s_{i,t}, a_{i,t+1}, Q_{i,t})\}$ 
12: Output:  $\hat{Q}(s, a)$ 

```

simulator is utilized [10] and N interactions are conducted between the simulated user and DM. Average cumulative reward is used as the fitness for the individual, which is similar to the objective of common RL algorithms.

$$F_R[\pi_{GA}] = \frac{1}{N} \sum_{i=1}^N \sum_{j=1}^{l_i} \gamma^{j-1} r_{ij} \quad (1)$$

where r_{ij} is the immediate reward and γ the discounted coefficient.

A noisy channel is designed to simulate ASR and SLU errors. For each dialog act $\{\text{act}, (\text{slot}, \text{value})\}$, replacement and deletion are randomly applied to *value* given the assigned confidence scores, which are randomly generated too. The produced N-best hypotheses are then fed into DMs.

2.3. On-corpus Q-points regression

Building a user simulator is not trivial and it is difficult to measure the consistency of the simulated user behavior to the real one. Learning a DM using a dialog corpus is appealing but there is very limited prior work on this subject [11, 12]. We propose to use an existing dialog corpus to estimate the fitness of a DM. First, an offline batch RL algorithm is applied on the corpus, inducing an optimum Q-function $\hat{Q}(s, a)$, and an implicitly defined policy $\pi_Q(s) = \arg \max_a \hat{Q}(s, a)$ which is optimum with respect to the corpus. Then $\hat{Q}(s, a)$ is used to define the fitness function. We use fitted Q-iteration [13] to learn a nonparametric approximator $\hat{Q}(s, a)$, as described in Algorithm 2. The algorithm uses Bellman equation (line 10) to update the estimated Q-values. Extremely Random Trees (ExtraTrees) [14] are utilized for non-parametric regression. ExtraTrees are a powerful model for regression and classification as they are both flexible and less susceptible to over-fitting. The annotated dialog corpus is represented as state-action-state triplets in the form of $\{(s_{t-1}, a_t, s_t)\}$, and

used as the training set. Two fitness estimation methods are proposed based on different heuristics. For an individual π_{GA} whose fitness to be evaluated, the NPoints fitness function is used to calculate the number of triplets where the actions predicted by π_{GA} and π_Q are identical.

$$F_{\text{NPoints}}[\pi_{GA}] = \sum_i \delta(\pi_{GA}(s_i), \pi_Q(s_i)) \quad (2)$$

The QVal fitness attempts to estimate the sum of the Q-values for the actions predicted by π_{GA} on the training triplets. However, the Q-function trained on a fixed corpus is often inaccurate in unexplored regions of the state space [15, 11]. To mitigate the problem a supervised classifier $\hat{P}(a|s)$ is built on the training set with the observed actions as targets. If the probability for an action is greater than a predefined threshold δ , the value produced by $\hat{Q}(s, a)$ is used, otherwise a constant R is used for punishment.

$$F_{\text{QVal}}[\pi_{GA}] = \sum_i \tilde{Q}_\delta(s_i, \pi_{GA}(s_i)) \quad (3)$$

$$\tilde{Q}_\delta(s, a) = \begin{cases} \hat{Q}(s, a) & \text{if } \hat{P}(a|s) > \delta \\ R & \text{otherwise} \end{cases}$$

The two fitness functions are different in weighing the importance of training instances. F_{QVal} will put a greater effort in optimizing instances with larger potential Q-value improvement while avoiding taking unobserved actions. Combining GA with the above two fitness functions leads to the on-corpus Q-points regression algorithm. One limitation of this algorithm compared to the on-line version is that no free parameter can be present in specifying the action structures since the fitness functions rely on the result of reinforcement learning, which does not support dynamical change of action structure.

2.4. On-corpus DM evaluation

We describe a DM evaluation method on dialog corpus without the need for deploying the DM online. A held-out dialog corpus is used as testing set, and the estimated cumulative reward for the testing dialogs when following the target DM policy is used as metric for performance. A similar approach has been taken in evaluating the effect of different dialog state tracker on end-to-end performance of a DM [15]. The estimation of Q-function is similar to Algorithm 2. But rather than learning the optimum policy, the value function for the policy to be evaluated is estimated, with the Bellman iteration (line 10) in Algorithm 2 changed to:

$$Q_{i,t} \leftarrow r + \gamma \hat{Q}(s_{i,t+1}, \pi(s_{i,t+1})) \quad (4)$$

where π is the DM policy to evaluate. Then the average reward for starting turns $\frac{1}{N} \sum_i Q_{i,0}$ is used as a metric for performance.

3. EXPERIMENTS

We devise a restaurant information domain for dialog simulation. There are 4 slots for the user simulator to fill before a database query. During simulations the DM interacts with the user simulator with a noise channel in between. The noise level of the channel can be adjusted to simulate different environmental noise conditions. Since the simulation process is stochastic, each experiment is conducted for 100 times, and the mean and standard deviation of testing performance are reported.

The reward function for the simulated environment is defined as follows. At each dialog turn the agent receives -1.0 reward. If correct restaurants are offered to users, 100.0 points are rewarded. But if the information is duplicate to that previously offered or the presented restaurants do not match user goal, -5.0 points are given. The reward discounting rate γ is set to 0.9.

In the on-corpus evaluation, DSTC2 dataset [16] is used for both DM policy learning and evaluation. The DSTC2 dataset was originally designed as a benchmark corpus for dialog state tracking. With the detailed annotation of dialog states, actions, SLU outputs and other information, it can be used as test set for end-to-end DM performance [15]. The dialog states used in both simulated and on-corpus experiments mainly comprise confident scores for each slot.

3.1. On-line learning experiment by simulation

The dialog policy template used in the simulated experiments is shown as follows.

- c0** On dialog beginning: `Welcome`
- c1** There are no valid SLU results or the top SLU hypothesis score is less than θ_0 : `Repeat`
- c2** User has just denied a slot: `Request that slot`
- c3** There is a slot with score less than θ_1 in the tracker:
if the score is larger than θ_2 then `ExplicitConf` else `Request`
- c4** The system has not yet output the action `RequireMore`:
`RequireMore`
- c5** Otherwise: query the database with slot-value pairs whose scores are greater than θ_3

The template has 6 condition-action clauses and contains 4 free parameters. Note that 3 free parameters lie in the condition expressions while θ_3 is used to adapt the semantics of the macro action `offer`, which queries the database and presents results to the user.

A rule-based DM policy is built by setting the 4 parameters heuristically. A RL-based policy trained using Q-learning with linear approximation is also built for comparison with the proposed methods. It is trained for over 100,000 dialog

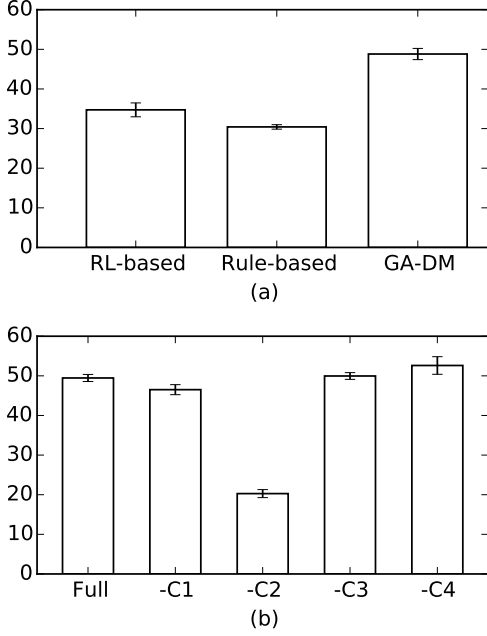


Fig. 2: Overall testing performance of DM policies. (a) Performance of RL-based, rule-based and GA-DM policy. (b) Performance of GA-DM with one clause disabled, where ‘-C1’ means that C1 is disabled.

sessions to ensure that the state-action space is sufficiently explored and the optimal performance is reached. The probability for exploration is set to 0.3. In the training process of each kind of DM, the noise level is randomly set for a dialog session. The same reward function and discounting rate are also used in the fitness estimation of GA. We run GA for 30 generations in policy training.

During testing 1000 dialog sessions are conducted and the noise level is adjusted in the same way as in training. We report both the overall performance (average reward received under a series of noise conditions) and the performance with a fixed noise level. The overall testing performance of each DM is shown in Fig.2. Performance when operating under fixed noise condition is shown in Fig.3 and 4. The level of environmental noise is measured using the semantic error rate of the top hypothesis of SLU results. It should be emphasized that the noise levels shown in the results are the same ones used in training. In addition to the GA-DM using the complete policy template, the utility of each individual clause in the template is evaluated. The four major clauses c1-c4 are disabled sequentially. The resulted DMs are evaluated using the same settings, and the testing results are shown along with the full-fledged GA-DM. The effects of different GA population size are explored and reported in Fig.5.

Since the DMs are optimized against the average reward received under several noise conditions, the overall testing reward shown in Fig.2 should be taken as the direct metric of performance. The RL-based policy showed better overall

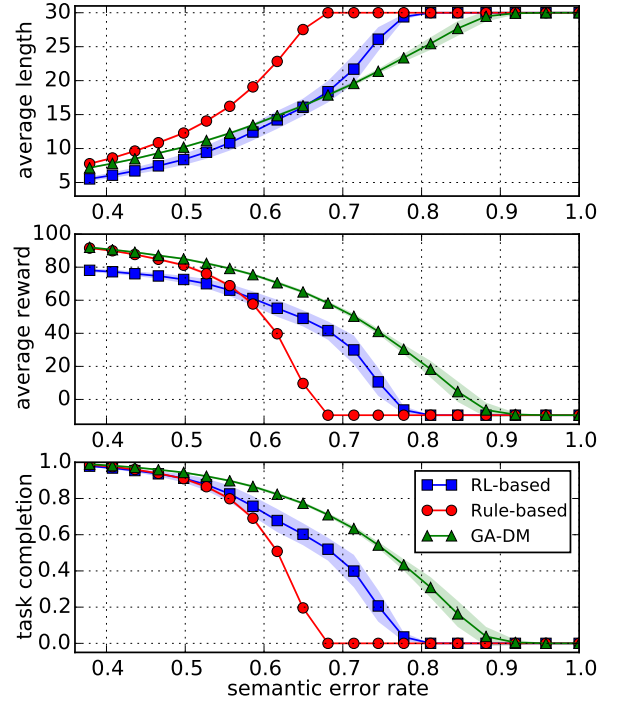


Fig. 3: On-line evaluation of GA-DM using a simulated user with fixed noise level. Average cumulative reward, dialog length and task completion rate are plotted against the error rate of the top SLU hypothesis, which is used as a metric for environmental noise.

performance than the rule-based one, while GA-DM significantly outperformed both the rule-based and RL-based policies. From Fig.3, it can be seen that when the noise is low, the rule-based DM is very competitive and shows even better performance than the RL-based DM. But when the noise level of the environment increases, its performance degrades seriously, while the RL-based DM is much more robust. However, after tuning of the free parameters using GA, the GA-DM outperforms both the other DMs on nearly all noise conditions. Note the maximum noise level at which each DM could successfully complete a dialog, suggesting that the GA-DM is able to operate under more adverse environment. It is worth mentioning again that the rule-based DM and GA-DM are instantiations of the same policy template. The simulation results justify GA as an effective method for DM policy optimization and reveal the performance potential of simple and yet human-interpretable DM policies.

It is interesting to make a comparison between RL and GA policy learning. In DM policy optimization, the state space is often continuous and infinite. In conventional RL, a model of the underlying optimal value function of the environment has to be designated. The ability of the model to approximate the optimal value function is a key factor affecting the performance of the learnt policy. However, the design of the model is often non-intuitive and complicated since it oper-

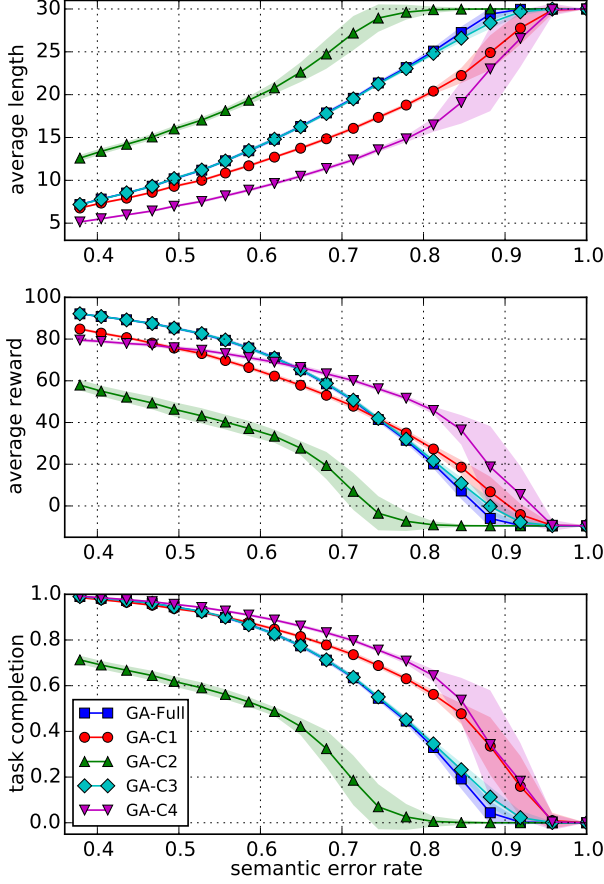


Fig. 4: Performance of the model operating with fixed noise level when the major clauses c1-c4 are sequentially disabled and retrained. Full GA-DM is the original model with all clauses enabled.

ates in the value function space. Expert knowledge is often difficult to be directly applied. This fact can help to explain that in our experiments, the RL-based DM is not as competitive as the others when the noise level is low. Since the noise level is varied during training, the resulted learning environment is much more difficult to deal with than one with fixed noise condition. Thus the linear model used is unlikely to perfectly match the underlying optimal value function and cannot accommodate all types of condition. In our experiment the RL policy has learnt to make a trade-off and adapted to conditions with high environmental noise for a better overall performance. GA-DM tackles the problem from a different perspective. It operates in policy space directly and is much easier to incorporate expert knowledge. In GA-DM a policy model is developed instead. Equivalent assumptions about policy structure are often difficult to made in value function space. Thus the resulted policy model can be more powerful and expressive than one for value function.

The relative utility of each clause of the policy template on the performance is another interesting aspect to be investi-

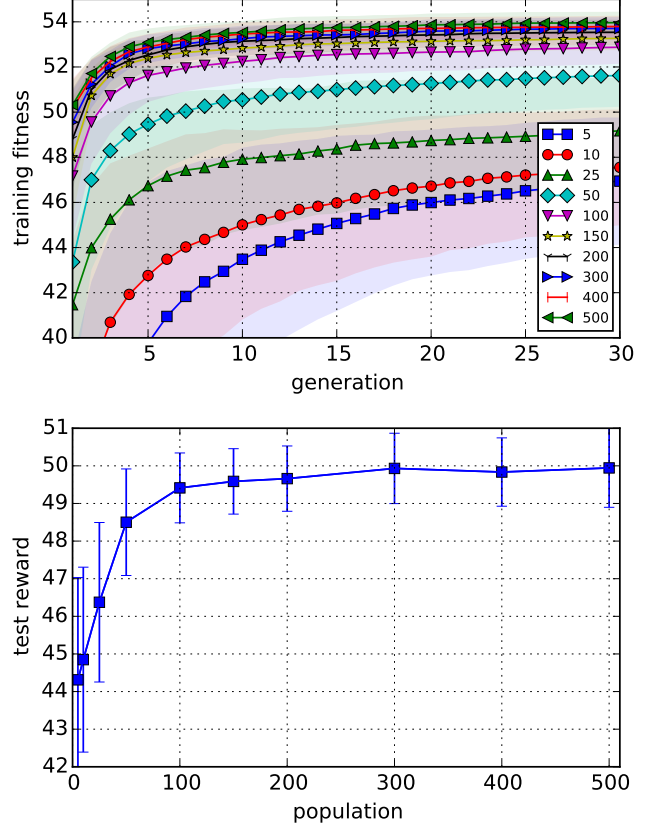


Fig. 5: Training fitness and testing performance of GA-DM when trained using different population size. Standard deviations are plotted as shaded areas and error bars.

gated. According to the results shown in Fig.2 (b) and Fig.4, it can be observed that when C2 is disabled the performance drops seriously. But to our surprise, when C4 is disabled, the performance significantly boosts especially in high-noise regions. The results show the relative utility of each clause in the template and reveal the necessity to optimize the structure of policy template. This kind of structural optimization problem can also be solved using GA, and we plan to study this kind of optimization in future work.

In GA the population size often influences the optimization efficiency. The training fitness and testing performance using different population size is shown in Fig.5. We can observe that with an increasing population size, the training and testing performance nearly monotonically increases. This performance improvements are more obvious when the size is less than 100, and are not noticeable above 300. Because the elitism technique is used and the fitness of the elitist individual is cached, the training fitness improves steadily during training.

3.2. On-corpus learning experiment

The DSTC2 testing corpus is used for on-corpus DM learning and evaluation [16], which is produced by a RL-based DM and consists of 1117 dialog sessions. The full annotations of the dataset are released after the conclusion of the DSTC2 challenge. The dialog state is the same as defined in the challenge, and we use the results produced by the ‘focus’ tracker using the scripts provided by the DSTC2 organizer. The dialog template used by GA comprises 9 condition-action clauses and 6 free parameters. The original corpus is equally split for training and testing.

The reward function is defined as follows. At each dialog turn the agent receives -10.0 reward. If correct restaurants are offered to users, 100.0 points are rewarded. But if the information is duplicate to that previously offered, -50.0 points are given. If the restaurants offered do not meet user’s demand, -100.0 points are given. The reward discounting rate γ is set to 0.9.

In addition to the GA-based DMs trained using QPoints-regression described in section 2.3, results of 3 additional DMs are shown for comparison.

1. SL-Original DM which is learned in a supervised way with the original dialog actions as training targets using the ExtraTrees classifier, represented as $\hat{P}(a|s)$.
2. SL-MaxQ supervised DM using the actions with maximum Q-value predicted by $\hat{Q}(s, a)$ as the supervised targets.
3. ThresholdedQ DM as described in [15], which selects the action with the maximum Q-value predicted by $\hat{Q}(s, a)$ from the set of actions whose probabilities produced by $\hat{P}(a|s)$ are greater than δ . The threshold is used to constrain the behavior of RL policy, in case of insufficient exploration.

To make full use of the available data and get a more stable estimation of the performance, we conducted 12 re-sampling experiments similar to the bootstrapping method, but avoid to use duplicate samples. In each sub-experiment, the dataset is reshuffled and split to get new training and testing instances. The average and standard deviation of the results are shown in Table 1.

The SL-MaxQ DM which acts greedily upon $\hat{Q}(s, a)$ has poor performance on the test set while being overrated on the training set, probably as a result of insufficient exploration. The ThresholdedQ DM mitigates the problem to a great degree by setting a simple threshold. That heuristic is shared with the QVal fitness function. GA-QVal outperforms all the other DMs and is very stable across the re-sampling experiments considering the relatively small standard deviation, while the behavior GA-NPoints which is less consistent results in an overall inferior performance. Although GA-QVal is trained under the guidance of an reinforcement learner $\hat{Q}(s, a)$, its performance is superior to both SL-MaxQ and

ThresholdedQ, which should be attributed to the prior domain knowledge incorporated into the policy template. The DMs in bold outperform SL-Original built by imitating the policy used in producing the corpus, indicating the possibility of building a better and yet human-comprehensible DM policy using a dialog corpus.

4. RELATED WORK

The subject of automatically optimizing dialog policies is a hot topic, and many data-driven methods have been proposed among which RL-based ones are the most popular. There is some previous work on constraining the behavior of RL-based DM. In [17] Williams proposed to construct a hand-crafted DM and it produces a set of candidate actions for given dialog state, from which the best one is chosen by a POMDP-DM. Lison [18] proposed to use ‘probabilistic rule’ in specifying the transition and reward sub-models of the POMDP model. The probabilistic rules are human-readable and less parameterized than conventional probability distribution, thus reducing the free parameters of the POMDP model and allowing the system designers to make use of domain knowledge in designing DM. Our work bears some resemblance to [18]. But we used the dialog policy template to specify a policy directly and utilized GA to train the free parameters.

Henderson et al. proposed a hybrid learning method in [11] to learn a policy on an existing dialog corpus by combining the results of supervised and reinforcement learning. Pure RL on fixed dataset often shows irregular behavior due to the insufficient exploration problem. Supervised learning (SL) is used to mitigate the problem and the hybrid method shows better performance than pure SL or RL. In this regard the QVal fitness function is similar in spirit and the use of policy template can further constrain the DM behavior, thus is suitable for off-line on-corpus learning.

One notable advantage of the GA-based DM over RL-based models is that the action structure can be changed during learning (only in on-line learning) as described in section 2.1. While in RL, each action $a_i \in A$ must be invariant otherwise the value function learned will be meaningless. This characteristic is suitable for SDS engineering since it can be difficult to determine the exact semantics of a dialog action beforehand. Further studies are needed in this regard.

5. CONCLUSIONS AND FUTURE WORK

In this paper we described a framework to train human-interpretable spoken dialog management policies using genetic algorithm. Two kinds of fitness functions were used, i.e., one based on interacting with a simulated user and the other on a dialog corpus which is more sample-efficient. We set up an online simulation environment and used the DSTC2 corpus for off-line on-corpus training and evaluation. The results show that by using domain language and setting appropriate

DM	Training	Testing
GA-NPoints	98.46 (38.30)	89.52 (41.30)
GA-QVal	127.38 (5.59)	129.29 (7.90)
SL-Original	115.63 (4.08)	114.39 (6.07)
SL-MaxQ	245.19 (12.59)	53.46 (36.06)
ThresholdedQ	142.48 (4.22)	122.21 (4.36)

Table 1: Estimated cumulative reward of DM policies on training and testing set. Numbers in brackets are standard deviation estimated by re-sampling experiments. Only starting turns of a dialog are considered as described in section 2.4. GA-NPoints and GA-QVal are DMs trained using GA with NPoints and QVal fitness functions respectively.

free parameters, the performance of simple rule-based DM policies can be largely improved, and can even outperform those trained using reinforcement learning. According to our knowledge, this is the first time that genetic algorithm is applied to DM optimization. Another advantage is its ability to optimize the structure of system actions. This framework is very suitable to upgrade existing SDSs using rule-based DM, by using collected data to optimize the newly specified free parameters.

This research is still preliminary and several aspects need further investigation, especially the effects of fitness functions. The search space of dialog policy in GA can be expanded by allowing the condition-action expressions to be reordered and partially disabled. The structural learning in system actions also needs further studies. We hope this work can help to build better and practical spoken dialog systems.

6. REFERENCES

- [1] Jason D. Williams and Steve Young, “Partially observable Markov decision processes for spoken dialog systems,” *Computer Speech & Language*, vol. 21, no. 2, pp. 393–422, 2007.
- [2] Steve Young, Milica Gašić, Simon Keizer, François Mairesse, Jost Schatzmann, Blaise Thomson, and Kai Yu, “The Hidden Information State model: A practical framework for POMDP-based spoken dialogue management,” *Computer Speech & Language*, vol. 24, no. 2, pp. 150–174, 2010.
- [3] Cheongjae Lee, Sangkeun Jung, Seokhwan Kim, and Gary Geunbae Lee, “Example-based dialog modeling for practical multi-domain dialog system,” *Speech Communication*, vol. 51, no. 5, pp. 466–484, 2009.
- [4] Pierre. Lison, *Structured Probabilistic Modelling for Dialogue Management*, Ph.D. thesis, University of Oslo, 2014.
- [5] Steve Young, Milica Gašić, Blaise Thomson, and Jason D. Williams, “POMDP-Based Statistical Spoken Dialog Systems: A Review,” *Proceedings of the IEEE*, vol. 101, no. 5, pp. 1160–1179, 2013.
- [6] Tim Paek, “Reinforcement Learning for Spoken Dialogue Systems: Comparing Strengths and Weaknesses for Practical Deployment,” Tech. Rep. MSR-TR-2006-62, Microsoft Research, 2006.
- [7] Yaakov Engel, Shie Mannor, and Ron Meir, “Reinforcement learning with Gaussian processes,” in *Proceedings of the 22nd international conference on Machine learning*. 2005, pp. 201–208, ACM.
- [8] Milica Gašić, Filip Jurčiček, Simon Keizer, François Mairesse, Blaise Thomson, Thomson, Kai Yu, and Steve Young, “Gaussian Processes for Fast Policy Optimisation of POMDP-based Dialogue Managers,” in *Proceedings of the 11th Annual Meeting of the Special Interest Group on Discourse and Dialogue*, Stroudsburg, PA, USA, 2010, SIGDIAL ’10, pp. 201–204, Association for Computational Linguistics.
- [9] Darrell Whitley, “A genetic algorithm tutorial,” *Statistics and Computing*, vol. 4, no. 2, pp. 65–85, June 1994.
- [10] Jost Schatzmann, Blaise Thomson, Karl Weilhammer, Hui Ye, and Steve Young, “Agenda-based user simulation for bootstrapping a POMDP dialogue system,” in *Human Language Technologies 2007: The Conference of the North American Chapter of the Association for Computational Linguistics; Companion Volume, Short Papers*. 2007, pp. 149–152, Association for Computational Linguistics.
- [11] James Henderson, Oliver Lemon, and Kallirroi Georgila, “Hybrid Reinforcement/Supervised Learning of Dialogue Policies from Fixed Data Sets,” *Computational Linguistics*, vol. 34, no. 4, pp. 487–511, 2008.
- [12] Olivier Pietquin, Matthieu Geist, Senthilkumar Chandramohan, and Hervé Frezza-Buet, “Sample-efficient batch reinforcement learning for dialogue management optimization,” *ACM Transactions on Speech and Language Processing (TSLP)*, vol. 7, no. 3, pp. 7, 2011.
- [13] Damien Ernst, Pierre Geurts, and Louis Wehenkel, “Tree-based batch mode reinforcement learning,” in *Journal of Machine Learning Research*, 2005, pp. 503–556.
- [14] Pierre Geurts, Damien Ernst, and Louis Wehenkel, “Extremely randomized trees,” *Machine Learning*, vol. 63, no. 1, pp. 3–42, Mar. 2006.

- [15] Sungjin Lee, “Extrinsic Evaluation of Dialog State Tracking and Predictive Metrics for Dialog Policy Optimization,” in *15th Annual Meeting of the Special Interest Group on Discourse and Dialogue*, 2014, p. 310.
- [16] Matthew Henderson, Blaise Thomson, and Jason D. Williams, “The second dialog state tracking challenge,” in *15th Annual Meeting of the Special Interest Group on Discourse and Dialogue*, 2014, p. 263.
- [17] Jason D. Williams, “The best of both worlds: unifying conventional dialog systems and POMDPs.,” in *INTER-SPEECH*, 2008, pp. 1173–1176.
- [18] Pierre Lison, “A hybrid approach to dialogue management based on probabilistic rules,” *Computer Speech & Language*, vol. 34, no. 1, pp. 232–255, 2015.